

Python で始める機械学習

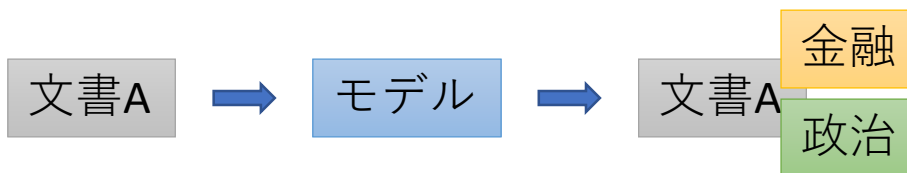
7.9 トピックモデリングと文書クラスタリング

15T4057F 藤井 真

7.9 トピックモデリングと文書クラスタリング

- トピックモデリング (topic modeling)

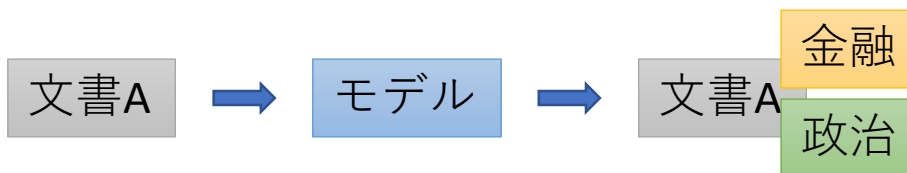
様々な文書に対して、1つ以上のtopicを割り当てるタスク。
通常、教師なし学習。



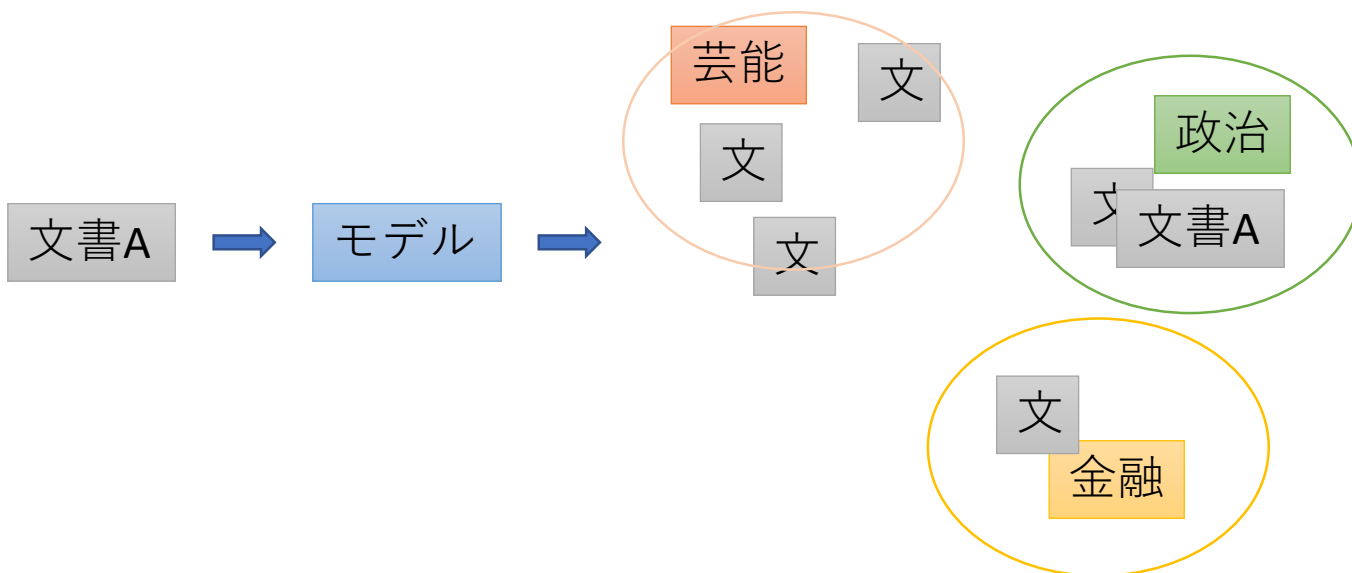
7.9 トピックモデリングと文書クラスタリング

- トピックモデリング (topic modeling)

様々な文書に対して、1つ以上のtopicを割り当てるタスク。
通常、教師なし学習。



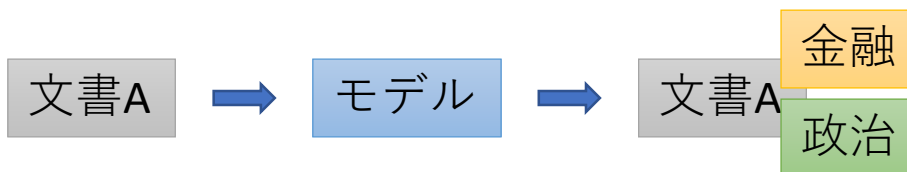
- 1文書に1トピックならクラスタリング的



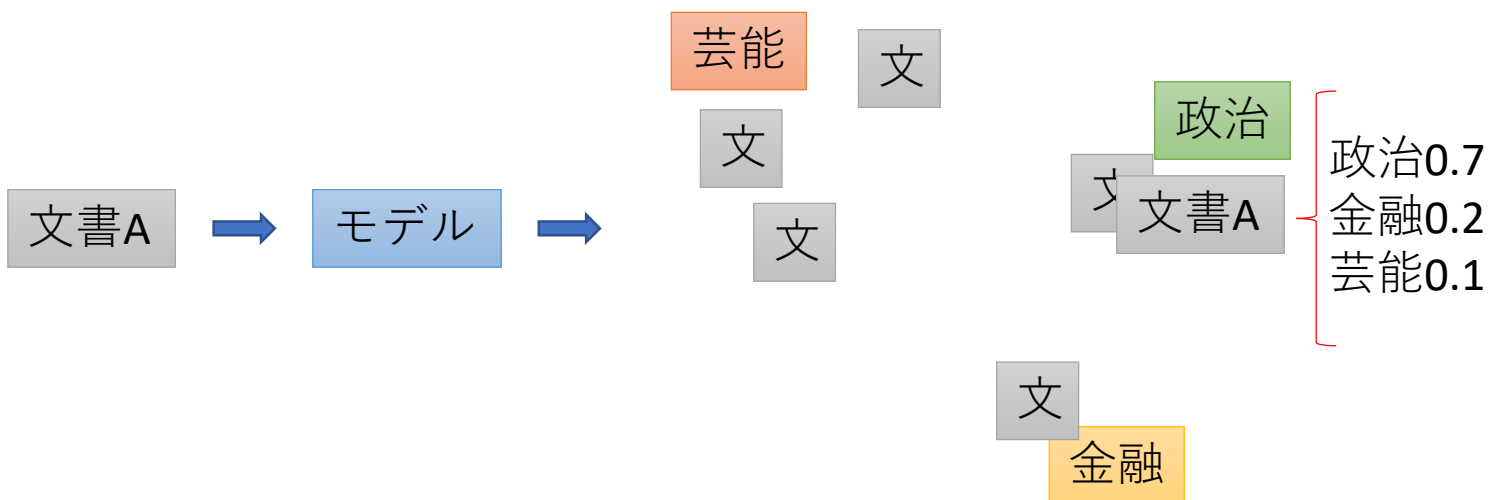
7.9 トピックモデリングと文書クラスタリング

- トピックモデリング (topic modeling)

様々な文書に対して、1つ以上のtopicを割り当てるタスク。
通常、教師なし学習。



- 1文書にトピックを2つ以上許せば成分分析的



7.9.1 LDA (Latent Dirichlet Allocation)

- LDA (Latent Dirichlet Allocation)
 - 隠れている、潜在的 ディリクレ分布 割り当て
 - トピックモデリングと言え、この手法と言える代表。
 - 文章は単語レベルの潜在的トピックの集まり。

例：BWEとWMDとWSDについて考えてください。・・・

7.9.1 LDA (Latent Dirichlet Allocation)

- LDA (Latent Dirichlet Allocation)

隠れている、潜在的 ディリクレ分布

割り当て

- トピックモデリングと言え、この手法と言える代表。
- 文章は単語レベルの潜在的トピックの集まり。

例：BWEとWMDとWSDについて考えてください。・・・

⇒各単語からこの文章トピックはNLP

- 逆に言えば同時に現れやすい単語の集合が文章になる。
- この同時に現れやすい単語を任意の数に分ければその数の分だけトピックに分けられたと解釈する。

7.9.1 LDA (Latent Dirichlet Allocation)

- LDAの映画レビューへの適用
 - データの準備

```
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.datasets import load_files

reviews_train = load_files("aclImdb/train/")
reviews_test = load_files("aclImdb/test/")
text_train, y_train = reviews_train.data, reviews_train.target
text_test, y_test = reviews_test.data, reviews_test.target

vect = CountVectorizer(max_features=10000, max_df=.15)
X = vect.fit_transform(text_train)
```

インポート

データロード

データ切り分け

BoW表現

- 文書の教師なし学習で、一般的な単語は学習を鈍らせることがあるのでBoWモデル作成時
 - 15%以上の文書に登場する単語を除く。
 - それ以外で最頻する10,000単語を用いる。

7.9.1 LDA (Latent Dirichlet Allocation)

- LDAの映画レビューへの適用
 - LDAモデルの作成

```
from sklearn.decomposition import LatentDirichletAllocation
                                10トピック
lda = LatentDirichletAllocation(n_topics=10, learning_method="batch",
                                max_iter=25, random_state=0)
document_topics = lda.fit_transform(X)                                結果の固定

print("lda.components_.shape: {}".format(lda.components_.shape))

lda.components_.shape: (10, 10000)
```

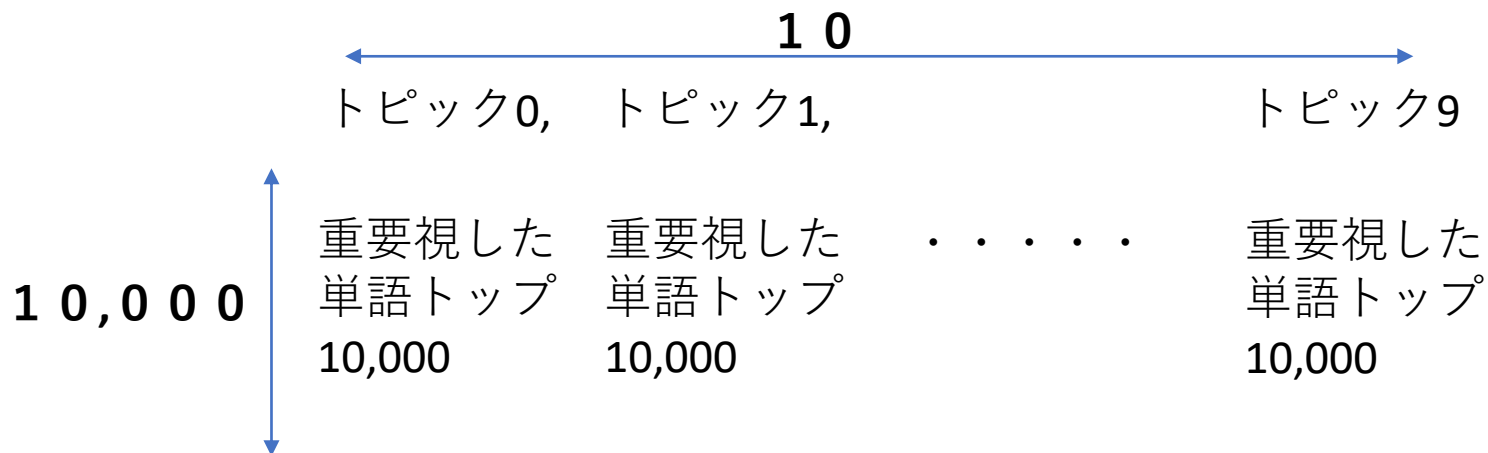
インポート

モデル作成

print文

出力

- 作成されたモデルイメージ



7.9.1 LDA (Latent Dirichlet Allocation)

- LDAの映画レビューへの適用
 - 10トピックの重要視した単語10個を見る。

```
import numpy as np
import mglearn

sorting = np.argsort(lda.components_, axis=1)[:, ::-1]

feature_names = np.array(vect.get_feature_names())

mglearn.tools.print_topics(topics=range(10), feature_names=feature_names,
                           sorting=sorting, topics_per_chunk=5, n_words=10)
```

10トピック

特徴量名称

ソート方法

10単語

インポート

ソート方法

単語名取得

表示設定

- print_topics関数を使うと結果をきれいに表示できる。

7.9.1 LDA (Latent Dirichlet Allocation)

- LDAの映画レビューへの適用
 - 10トピックの重要視した単語10個を見る。

topic 0	topic 1	topic 2	家系?	コメディ?
-----	-----	-----	-----	-----
director	show	book	family	funny
work	series	original	young	comedy
performance	war	10	father	cast
actors	episode	now	us	role
cast	tv	again	woman	humor
screen	years	world	own	fun
performances	american	saw	world	jokes
role	episodes	read	real	actors
both	world	didn	mother	performance
quite	shows	am	between	always

ゾンビホラー?

topic 5	topic 6	topic 7	topic 8	警察系アクション?
-----	-----	-----	-----	-----
horror	music	original	thing	action
gore	john	team	worst	police
effects	old	series	didn	murder
blood	young	jack	nothing	killer
pretty	girl	action	minutes	crime
budget	song	new	guy	plays
house	gets	down	actually	lee
zombie	dance	tarzan	want	gets
dead	songs	freddy	going	role
low	rock	indian	re	cop

7.9.1 LDA (Latent Dirichlet Allocation)

- LDAの映画レビューへの適用
 - 10トピックの重要視した単語10個を見る。

topic 0 -----	topic 1 -----	topic 2 -----	家系? topic 3 -----	コメディ? topic 4 -----
director	show	book	family	funny
work	series	original	young	comedy
performance	war	10	father	cast
actors	episode	now	us	role
cast	tv	again	woman	humor
screen	years	world	own	fun
performances	american	saw	world	jokes
role	episodes	read	real	actors
both	world	didn	mother	performance
quite	shows	am	between	always

ゾンビホラー?

topic 5 -----	topic 6 -----	topic 7 -----	topic 8 -----	警察系アクション? topic 9 -----
horror	music	original	thing	action
gore	john	team	worst	police
effects	old	series	didn	murder
blood	young	jack	nothing	killer
pretty	girl	action	minutes	crime
budget	song	new	guy	plays
house	gets	down	actually	lee
zombie	dance	tarzan	want	gets
dead	songs	freddy	going	role
low	rock	indian	re	cop



10トピックで映画全体をカバーするのは、1トピックあたりが幅広くなってぼやける。

7.9.1 LDA (Latent Dirichlet Allocation)

- LDAの映画レビューへの適用
 - 100トピックなら？

```
lda100 = LatentDirichletAllocation(n_topics=100, learning_method="batch",  
                                  max_iter=25, random_state=0)  
document_topics100 = lda100.fit_transform(X)  
  
sorting = np.argsort(lda100.components_, axis=1)[: , :-1]  
feature_names = np.array(vect.get_feature_names())  
topics_chociced = np.array([3,4,5,9,23,56,76,78,79,81])  
print("\n")  
mglearn.tools.print_topics(topics=topics_chociced, feature_names=feature_names,  
                           sorting=sorting, topics_per_chunk=5, n_words=10)
```

- 10トピックとほぼ同様。
- 100のうち分かりやすい10トピックを選択。

7.9.1 LDA (Latent Dirichlet Allocation)

- LDAの映画レビューへの適用
 - 100トピックなら？

キリスト教

topic 3

god
church
christian
jesus
religious
faith
religion
message
christ
believe

バットマン

topic 4

batman
jane
welles
mr
rochester
eyre
timothy
kane
pitt
robin

ゾンビとマイケル ヴァンプ ジャクソン

topic 5

horror
zombie
gore
zombies
blood
vampire
dead
vampires
pretty
house

topic 9

music
michael
band
voice
jackson
video
rock
price
fan
clips

日本の アニメ

topic 23

japanese
animation
quite
animated
japan
anime
world
bit
however
english

自分の興味ある
トピックの単語を
確認してみてください。

トムと ジェリー

topic 56

tom
jerry
russell
may
anne
kurt
mouse
andrew
night
darren

猟奇

topic 76

killer
girl
killed
death
slasher
killing
kill
gets
horror
dead

家系

topic 78

mother
family
child
daughter
father
son
husband
real
boy
heart

007

topic 79

peter
allen
woody
bond
james
brozman
falk
sellers
role
bergman

STAR WARS

topic 81

star
wars
luke
washington
joseph
battle
lucas
smith
denzel
fisher

7.9.1 LDA (Latent Dirichlet Allocation)

- LDAの映画レビューへの適用
 - 100トピックなら？

キリスト教

topic 3

god
church
christian
jesus
religious
faith
religion
message
christ
believe

バットマン

topic 4

batman
jane
welles
mr
rochester
eyre
timothy
kane
pitt
robin

ゾンビとマイケル ヴァンプ ジャクソン

topic 5

horror
zombie
gore
zombies
blood
vampire
dead
vampires
pretty
house

topic 9

music
michael
band
voice
jackson
video
rock
price
fan
clips

日本の アニメ

topic 23

japanese
animation
quite
animated
japan
anime
world
bit
however
english

自分の興味ある
トピックの単語を
確認してみてください。

トムと ジェリー

topic 56

tom
jerry
russell
may
anne
kurt
mouse
andrew
night
darren

猟奇

topic 76

killer
girl
killed
death
slasher
killing
kill
gets
horror
dead

家系

topic 78

mother
family
child
daughter
father
son
husband
real
boy
heart

007

topic 79

peter
allen
woody
bond
james
brozman
falk
sellers
role
bergman

STAR WARS

topic 81

star
wars
luke
washington
joseph
battle
lucas
smith
denzel
fisher



10トピックより具体的

7.9.1 LDA (Latent Dirichlet Allocation)

- LDAの映画レビューへの適用

- トピック 81 : STARWARS のレビューを確認する。

```
starwars = np.argsort(document_topics100[:, 81])[:-1]
for i in starwars[:5]: 5文書
    print(b".".join(text_train[i].split(b" ")[[:2]] + b".\n"))
```

○ b'Despite the feelings of most "Star Wars" fans, in my opinion "Return Of The Jedi" is the gr
eatest cinematic film ever created. Ever since the first time I saw it, it\'s depth, intensit
y, special effects, and moving story have overwhelmed me.\n'

△
✕ b'A very good movie. A classic sci-fi film with humor, action and everything.\n'

b"He glorified himself as a great supporting actor in `Glory`, he proved he was no `Malcolm i
n the Middle` mediocre actor in `Malcolm X`, he showed his brotherly love for acting in `Phil
adelphia`, he pulled a slam dunk in `He Got Game`, he pulled no punches and rocked us like a
hurricane in `The Hurricane`, he provided us effective thespian education in `Training Day`,
and now he has demonstrated that he could also direct! Denzel Washington's directorial debut
`Antwone Fisher` is the most moving film of the year. This tearjerker `fish`er story is in no
relation to the debacle that happened to the Miami Dolphins in the 4th quarter against the Ne
w England Patriots in the last game of the 2002 season.\n"

○ b'Return of the Jedi is certainly the most action packed of the series, and is a fine conclus
ion to the Star Wars Saga. With Han Solo imprisoned by Jabba the Hut and the Empire building
a new Death Star, the rebel alliance is facing an uphill struggle against the dark side, and
only our favourite heroes can pull it off.\n'

○ b"I noticed that A NEW HOPE and THE EMPIRE STRIKES BACK are in the TOP 10, but that this one
isn't even in the TOP 100.

This movie has a bad reputation because of Ewoks, but t
here are so many reasons to love this movie:

-The Rescue of Han Solo from Jabba: T
his official wraps up the Han Solo in debt sub-plot that was established when we first met th
e character in A NEW HOPE.\n"

3文書目は1文あたりが長く、特殊な内容⇒特徴がとらえづらい

7.9.1 LDA (Latent Dirichlet Allocation)

- LDAの映画レビューへの適用
 - 100トピックが全文書に対して得た重み

```
import matplotlib.pyplot as plt

fig, ax = plt.subplots(1, 2, figsize=(10, 10))
topic_names = ["{:>2} ".format(i) + " ".join(words)
               for i, words in enumerate(feature_names[sorting[:, :2]])]

for col in [0, 1]:
    start = col * 50
    end = (col + 1) * 50
    ax[col].barh(np.arange(50), np.sum(document_topics100, axis=0)[start:end])
    ax[col].set_yticks(np.arange(50))
    ax[col].set_yticklabels(topic_names[start:end], ha="left", va="top")
    ax[col].invert_yaxis()
    ax[col].set_xlim(0, 2000)
    yax = ax[col].get_yaxis()
    yax.set_tick_params(pad=130)

plt.tight_layout()
plt.show()
```

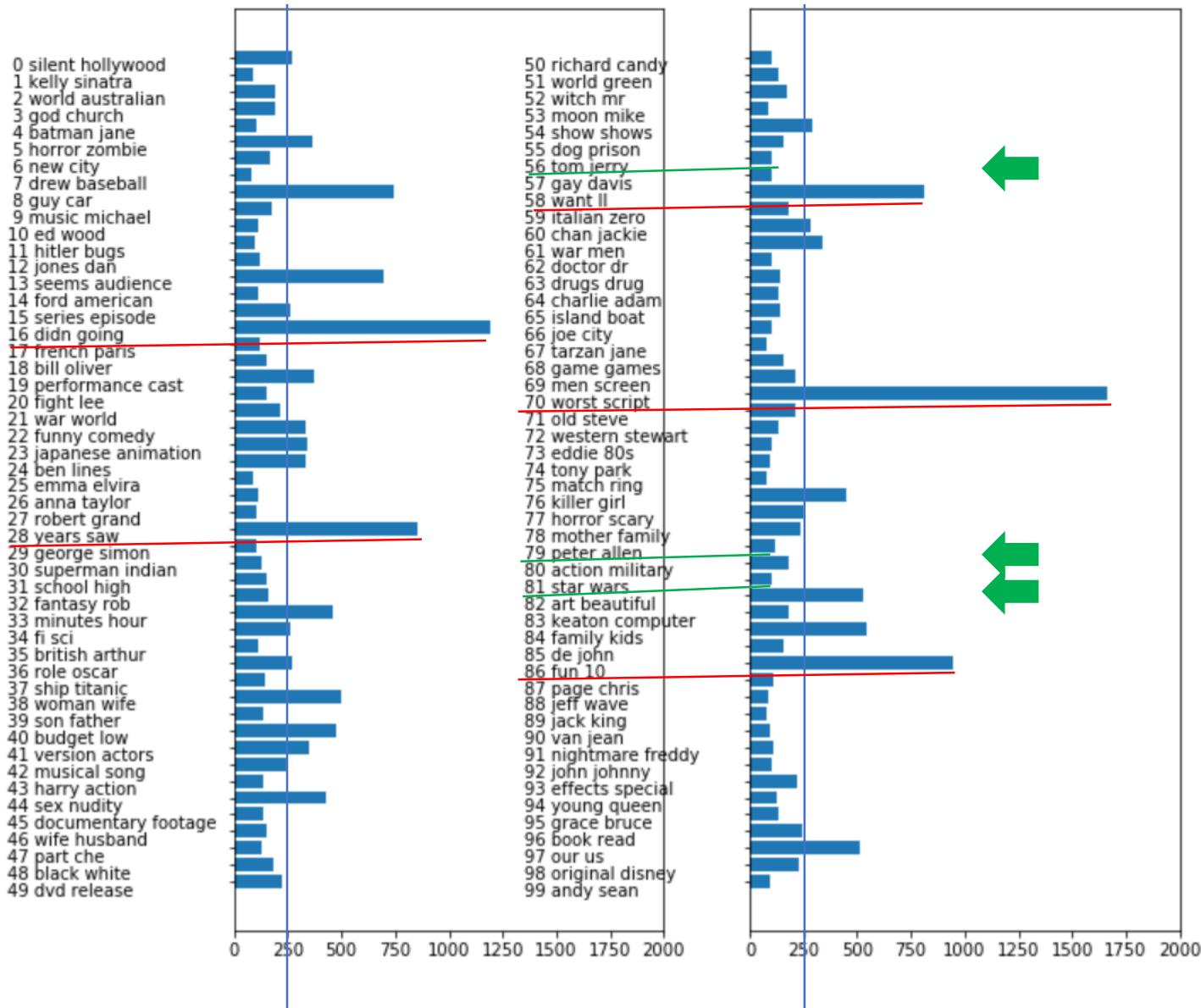
代表的な2単語をトピックタイトルとする

インポート

目盛り等
設定

2列で表示
値の取得

7.9.1 LDA (Latent Dirichlet Allocation)



- 25,000 文書
100 トピック
平均は 250 文書
- 赤線は抽象的な好悪文
文章を集めやすい
- 緑線は上から
 - トムとジェリー
 - 007
 - STARWARS
 具体的で文章が
集まらない

7.9.1 LDA (Latent Dirichlet Allocation)

- まとめ

- LDAは「潜在的ディリクレ(分布)割り当て」
- トピックモデリングを代表する教師なし学習手法
- トピック数を増やすと具体的トピックが得られる
- 教師なしらしく人の解釈がいるため
ただトピック数を増やせばよいわけではない